# Operator Language: A Program Generation Framework for Fast Kernels

**Franz Franchetti**, **Frédéric de Mesmay, Daniel McFarlin, Markus Püschel**
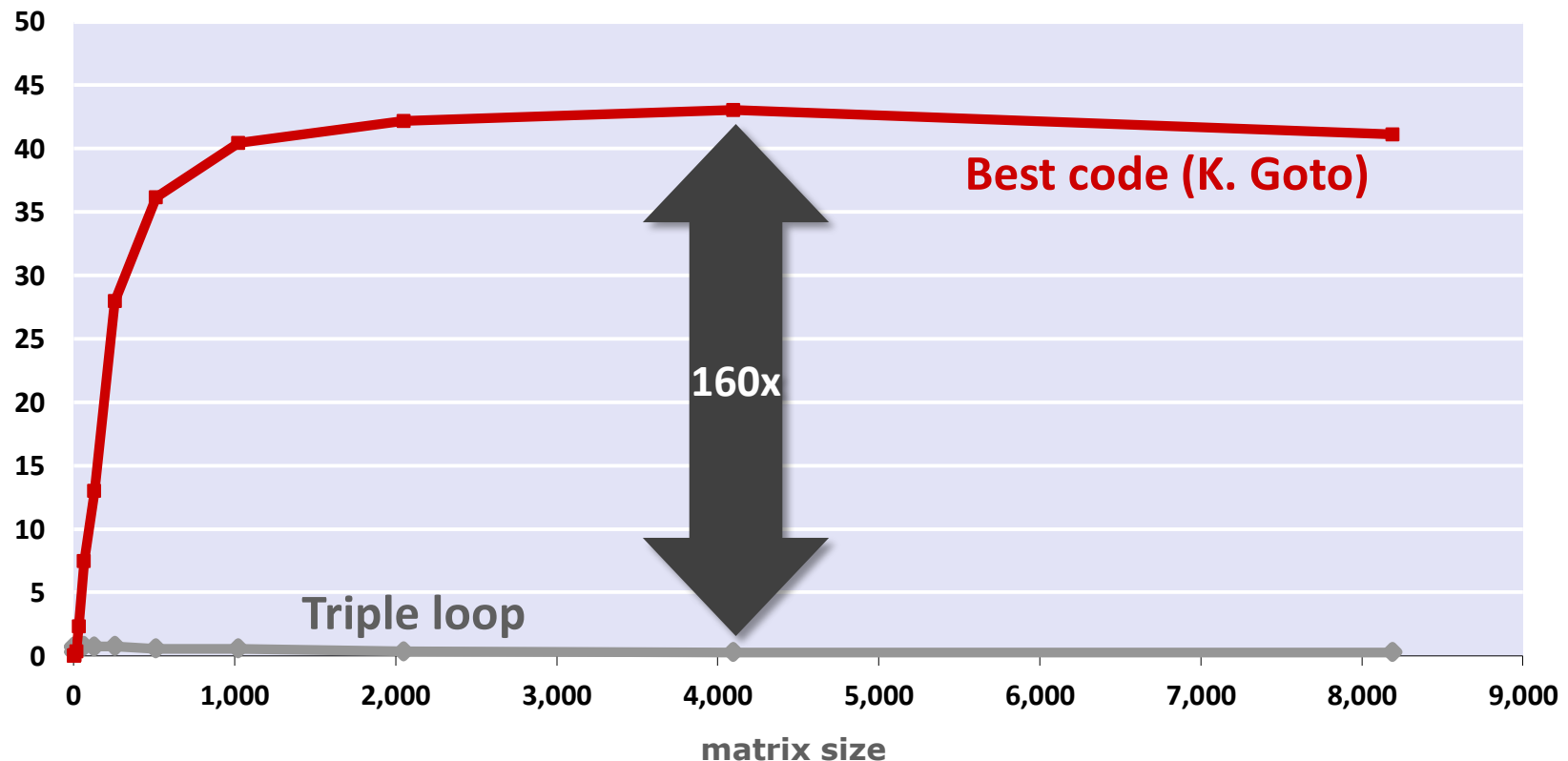
**Electrical and Computer Engineering**
**Carnegie Mellon University**

# The Problem: Example MMM

**Matrix-Matrix Multiplication (MMM) on 2xCore2Duo 3 GHz (double precision)**
**Performance [Gflop/s]**



**Best code (K. Goto)**

**160x**

**Triple loop**

matrix size

- **Similar plots can be shown for all numerical kernels in linear algebra, signal processing, coding, crypto, …**

- *What's going on? Hardware is becoming increasingly complex.*
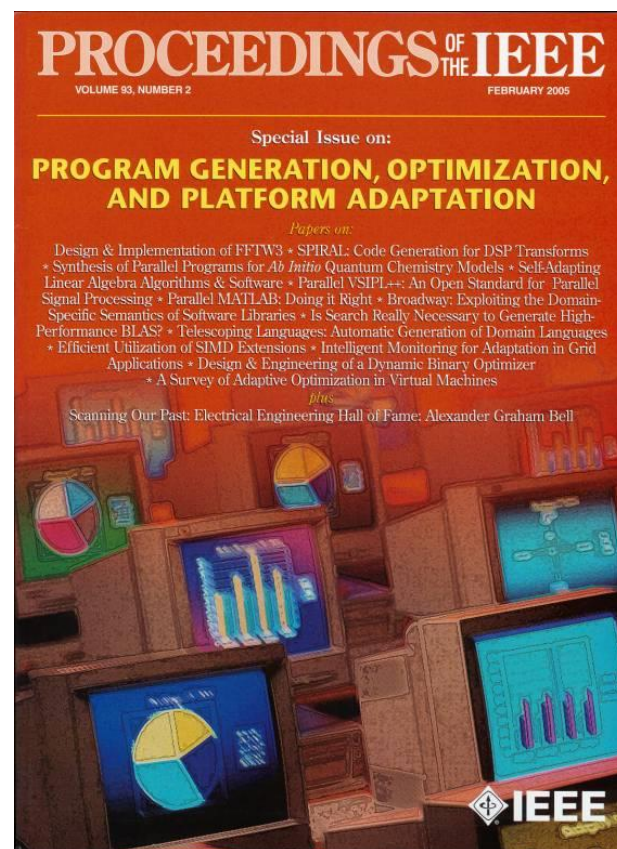
# Automatic Performance Tuning

- **Current vicious circle: Whenever a new platform comes out, the same functionality needs to be rewritten and reoptimized**

- **Automatic Performance Tuning**
  - BLAS: ATLAS, PHiPAC
  - Linear algebra: Sparsity/OSKI, Flame
  - Sorting
  - Fourier transform: FFTW
  - Linear transforms (and beyond): Spiral
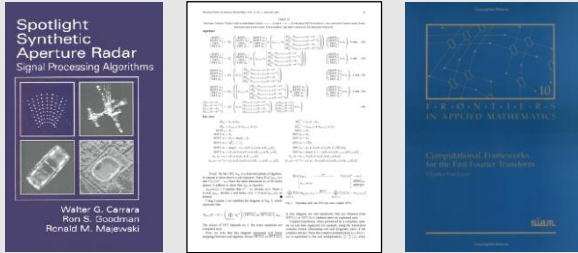  - …others

  *How to build an extensible system?*
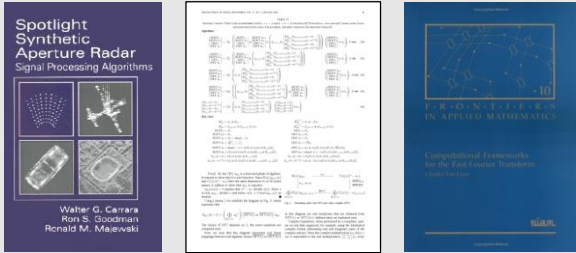  *For more problem classes?*
  *For yet un-invented platforms?*



**PROCEEDINGS OF THE IEEE**
VOLUME 93, NUMBER 2   FEBRUARY 2005

Special Issue on:
**PROGRAM GENERATION, OPTIMIZATION, AND PLATFORM ADAPTATION**

*Papers on:*
Design & Implementation of FFTW3 • SPIRAL: Code Generation for DSP Transforms • Synthesis of Parallel Programs for *Ab Initio* Quantum Chemistry Models • Self-Adapting Linear Algebra Algorithms & Software • Parallel VSIPL++: An Open Standard for Parallel Signal Processing • Parallel MATLAB: Doing it Right • Broadway: Exploiting the Domain-Specific Semantics of Software Libraries • Is Search Really Necessary to Generate High-Performance BLAS? • Telescoping Languages: Automatic Generation of Domain Languages • Efficient Utilization of SIMD Extensions • Intelligent Monitoring for Adaptation in Grid Applications • Design & Engineering of a Dynamic Binary Optimizer • A Survey of Adaptive Optimization in Virtual Machines
*plus*
Scanning Our Past: Electrical Engineering Hall of Fame: Alexander Graham Bell

**Proceedings of the IEEE special issue, Feb. 2005**

# What is Spiral?



*Traditionally*

*Spiral Approach*

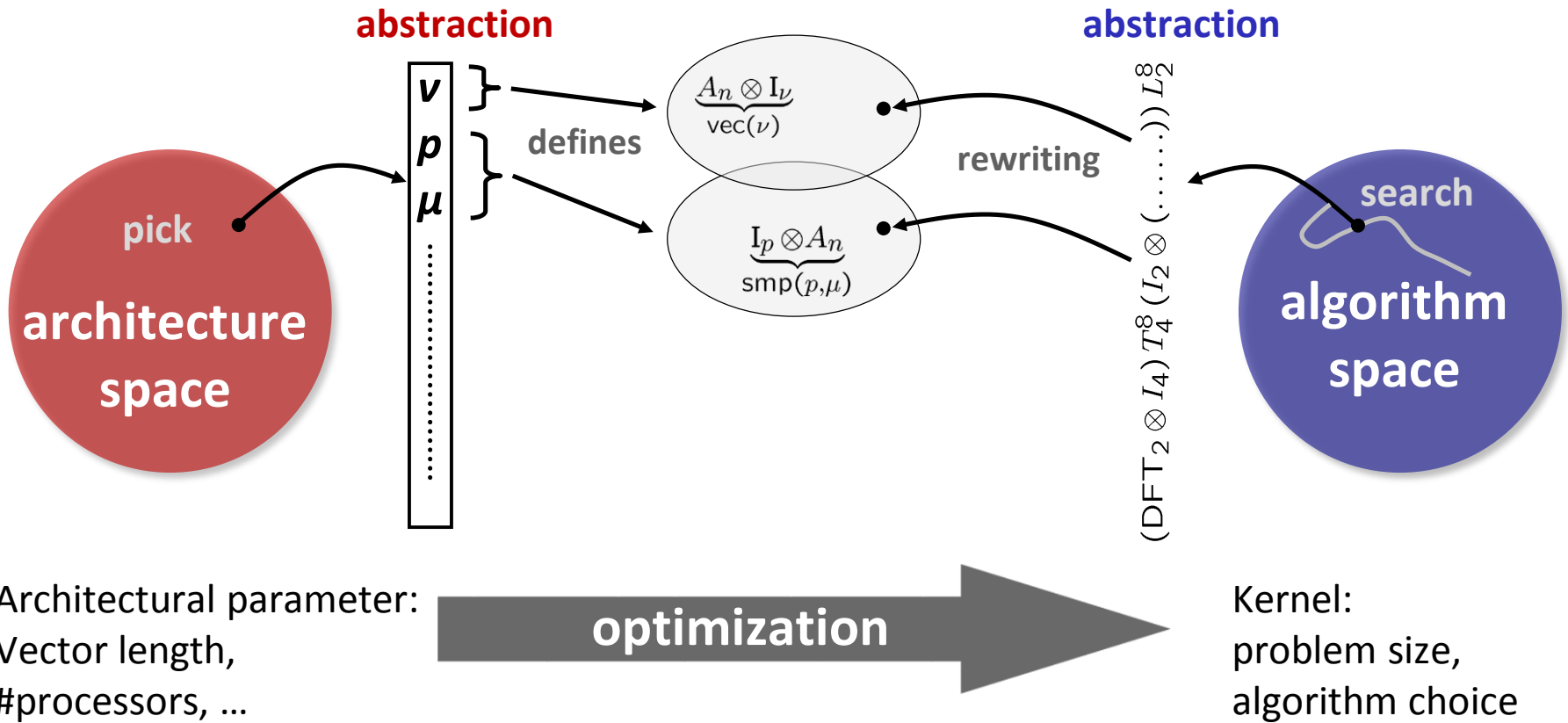**High performance library optimized for given platform**

*Comparable performance*

**High performance library optimized for given platform**

# Idea: Common Abstraction and Rewriting

**Model:** common abstraction
= spaces of matching formulas
= domain-specific language



**abstraction**

**abstraction**

$\dfrac{A_n \otimes I_\nu}{\mathrm{vec}(\nu)}$

$\dfrac{I_p \otimes A_n}{\mathrm{smp}(p,\mu)}$

$\nu$
$p$
$\mu$

**defines**

**rewriting**

**pick**

**search**

**architecture space**

**algorithm space**

$(\mathrm{DFT}_2 \otimes I_4)\, T_4^8\, (I_2 \otimes (\ldots\ldots))\, L_2^8$

Architectural parameter:
Vector length,
#processors, …

**optimization**

Kernel:
problem size,
algorithm choice

# Some Kernels as OL Formulas.

## Linear Transforms

$$\mathbf{DFT}_n \;\rightarrow\; (\mathbf{DFT}_k \otimes \mathrm{I}_m)\, \mathsf{T}^n_m (\mathrm{I}_k \otimes \mathbf{DFT}_m)\, \mathsf{L}^n_k, \quad n = km$$

$$\mathbf{DFT}_n \;\rightarrow\; P_n(\mathbf{DFT}_k \otimes \mathbf{DFT}_m)Q_n, \quad n = km,\ \gcd(k,m) = 1$$

$$\mathbf{DFT}_p \;\rightarrow\; R_p^T(\mathrm{I}_1 \oplus \mathbf{DFT}_{p-1})D_p(\mathrm{I}_1 \oplus \mathbf{DFT}_{p-1})R_p, \quad p \text{ prime}$$

$$\mathbf{DCT\text{-}3}_n \;\rightarrow\; (\mathrm{I}_m \oplus \mathsf{J}_m)\, \mathsf{L}^n_m (\mathbf{DCT\text{-}3}_m(1/4) \oplus \mathbf{DCT\text{-}3}_m(3/4))$$

$$\cdot (\mathsf{F}_2 \otimes \mathrm{I}_m) \begin{bmatrix} \mathrm{I}_m & 0 \oplus -\mathsf{J}_{m-1} \\ & \frac{1}{\sqrt{2}}(\mathrm{I}_1 \oplus 2\,\mathrm{I}_m) \end{bmatrix}, \quad n = 2m$$

$$\mathbf{DCT\text{-}4}_n \;\rightarrow\; S_n \mathbf{DCT\text{-}2}_n \operatorname{diag}_{0 \le k < n}(1/(2\cos((2k+1)\pi/4n)))$$

$$\mathbf{IMDCT}_{2m} \;\rightarrow\; (\mathsf{J}_m \oplus \mathrm{I}_m \oplus \mathrm{I}_m \oplus \mathsf{J}_m)\left(\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix} \otimes \mathrm{I}_m\right) \oplus \left(\begin{bmatrix} -1 \\ -1 \end{bmatrix} \otimes \mathrm{I}_m\right)\right) \mathsf{J}_{2m} \mathbf{DCT\text{-}4}_{2m}$$

$$\mathbf{WHT}_{2^k} \;\rightarrow\; \prod_{i=1}^{t} (\mathrm{I}_{2^{k_1 + \cdots + k_{i-1}}} \otimes \mathbf{WHT}_{2^{k_i}} \otimes \mathrm{I}_{2^{k_{i+1} + \cdots + k_t}}), \quad k = k_1 + \cdots + k_t$$

$$\mathbf{DFT}_2 \;\rightarrow\; \mathsf{F}_2$$

$$\mathbf{DCT\text{-}2}_2 \;\rightarrow\; \operatorname{diag}(1, 1/\sqrt{2})\, \mathsf{F}_2$$

$$\mathbf{DCT\text{-}4}_2 \;\rightarrow\; \mathsf{J}_2\, \mathsf{R}_{13\pi/8}$$

## Viterbi Decoding



010001 convolutional encoder  11 10 00 01 10 01 11 00  11 10 01 01 10 10 11 00  Viterbi decoder  010001

$$\underbrace{\mathbf{Vit}}_{\mathrm{vec}(v)} \;\rightarrow\; \underbrace{\left(\prod (L \times I) \circ (I \otimes C)\right) \circ Id}_{\mathrm{vec}(v)}$$

$$\rightarrow\; \left(\prod \underbrace{(L \times I) \circ (I \otimes C)}_{\mathrm{vec}(v)}\right) \circ Id$$

$$\pounds \quad \rightarrow\; \left(\prod (L \otimes I_v \times I) \circ (I \otimes C \otimes I_v) \circ (\vec{L} \times I)\right) \circ Id$$

$$\rightarrow\; \prod (L \otimes I_v \times I) \circ (I \otimes (B \otimes I_v)) \circ (\vec{L} \times I)$$

## Matrix-Matrix Multiplication



$$= \quad \pounds$$

$$\mathrm{MMM}_{1,1,1} \rightarrow (\cdot)_1$$

$$\mathrm{MMM}_{m,n,k} \rightarrow (\otimes)_{m/m_b \times 1} \otimes \mathrm{MMM}_{m_b,n,k}$$

$$\mathrm{MMM}_{m,n,k} \rightarrow \mathrm{MMM}_{m,nb,k} \otimes (\otimes)_{1 \times n/nb}$$

$$\mathrm{MMM}_{m,n,k} \rightarrow ((\Sigma_{k/k_b} \circ (\cdot)_{k/k_b}) \otimes \mathrm{MMM}_{m,n,k_b}) \circ$$
$$((L^{mk/k_b}_{k/k_b} \otimes I_{k_b}) \times I_{kn})$$

$$\mathrm{MMM}_{m,n,k} \rightarrow (L^{mn/n_b}_m \otimes I_{n_b}) \circ$$
$$((\otimes)_{1 \times n/n_b} \otimes \mathrm{MMM}_{m,n_b,k}) \circ$$
$$(I_{km} \times (L^{kn/n_b}_{n/n_b} \otimes I_{n_b}))$$

## Synthetic Aperture Radar (SAR)



preprocessing → matched filtering → interpolation → 2D iFFT

$$\mathsf{SAR}_{k \times m \to n \times n} \;\rightarrow\; \mathsf{DFT}_{n \times n} \circ \mathsf{Interp}_{k \times m \to n \times n}$$

$$\mathsf{DFT}_{n \times n} \;\rightarrow\; (\mathsf{DFT}_n \otimes \mathrm{I}_n) \circ (\mathrm{I}_n \otimes \mathsf{DFT}_n)$$

$$\mathsf{Interp}_{k \times m \to n \times n} \;\rightarrow\; (\mathsf{Interp}_{k \to n} \otimes_i \mathrm{I}_n) \circ (\mathrm{I}_k \otimes_i \mathsf{Interp}_{m \to n})$$

$$\mathsf{Interp}_{r \to s} \;\rightarrow\; \left(\bigoplus_{i=0}^{n-2} \mathsf{InterpSeg}_k\right) \oplus \mathsf{InterpSegPruned}_{k,\ell}$$

$$\mathsf{InterpSeg}_k \;\rightarrow\; \mathsf{G}^{u \cdot n \to k}_f \circ \mathsf{iPrunedDFT}_{n \to u \cdot n} \circ \left(\frac{1}{n}\right) \circ \mathsf{DFT}_n$$

# How Spiral Works

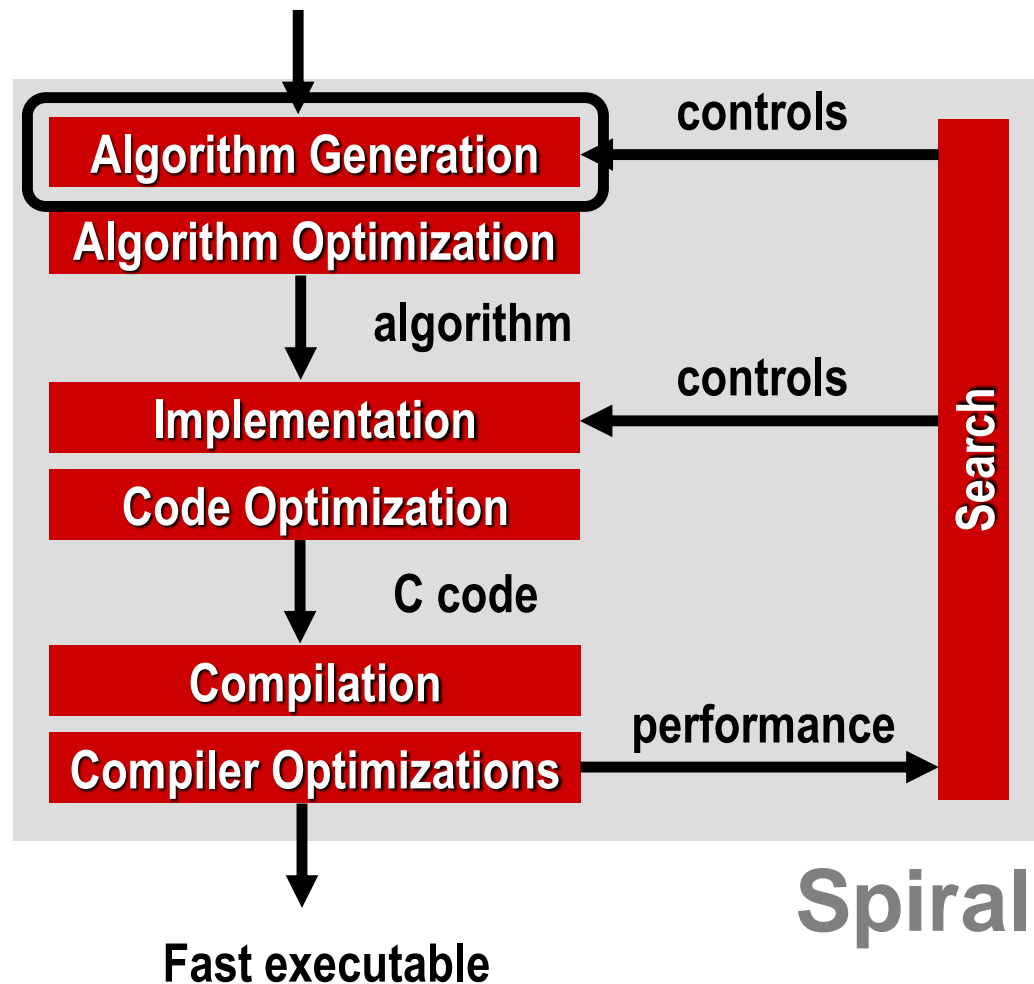**Problem specification (transform)**

**Spiral:**

**Complete automation of the implementation and optimization task**

**Basic ideas:**

**Declarative representation of algorithms**

**Rewriting systems to generate and optimize algorithms at a high level of abstraction**



Algorithm Generation
Algorithm Optimization
controls
algorithm
Implementation
controls
Code Optimization
C code
Compilation
Compiler Optimizations
performance
Search
Spiral

**Fast executable**

Markus Püschel, José M. F. Moura, Jeremy Johnson, David Padua, Manuela Veloso, Bryan Singer, Jianxin Xiong, Franz Franchetti, Aca Gacic, Yevgen Voronenko, Kang Chen, Robert W. Johnson, and Nick Rizzolo:

# Organization

- **Operator language and algorithms**

- **Optimizing algorithms for platforms**

- **Performance results**

- **Summary**

# Organization

- **Operator language and algorithms**

- Optimizing algorithms for platforms

- Performance results

- Summary

# Operators

## Definition

- **Operator: Multiple complex vectors ! multiple complex vectors**
- **Higher-dimensional data is linearized**
- **Operators are potentially nonlinear**

$$\mathsf{M} : \begin{cases} \mathbb{C}^{n_0} \times \cdots \times \mathbb{C}^{n_{k-1}} \longrightarrow \mathbb{C}^{N_0} \times \cdots \times \mathbb{C}^{N_{\ell-1}} \\ (\mathbf{x}_0, \mathbf{x}_1, \ldots, \mathbf{x}_{k-1}) \mapsto \mathsf{M}(\mathbf{x}_0, \mathbf{x}_1, \ldots, \mathbf{x}_{k-1}) \end{cases}$$

## Example: Matrix-matrix-multiplication (MMM)

$$\mathrm{MMM}_{m,k,n} : \mathbb{R}^{mk} \times \mathbb{R}^{kn} \rightarrow \mathbb{R}^{mn}$$
$$(\mathbf{A}, \mathbf{B}) \mapsto \mathbf{AB}$$

# Operator Language

| name | definition |
|------|-----------|
| *Linear, arity (1,1)* | |
| identity | $I_n : \mathbb{C}^n \to \mathbb{C}^n; \mathbf{x} \mapsto \mathbf{x}$ |
| vector flip | $J_n : \mathbb{C}^n \to \mathbb{C}^n; (x_i) \mapsto (x_{n-i})$ |
| transposition of an $m \times n$ matrix | $L_m^{mn} : \mathbb{C}^{mn} \to \mathbb{C}^{mn}; \mathbf{A} \mapsto \mathbf{A}^T$ |
| matrix $M \in \mathbb{C}^{m \times n}$ | $\mathrm{M} : \mathbb{C}^n \to \mathbb{C}^m; \mathbf{x} \mapsto M\mathbf{x}$ |
| *Multilinear, arity (2,1)* | |
| Point-wise product | $\mathrm{P}_n : \mathbb{C}^n \times \mathbb{C}^n \to \mathbb{C}^n; ((x_i),(y_i)) \mapsto (x_i y_i)$ |
| Scalar product | $\mathrm{S}_n : \mathbb{C}^n \times \mathbb{C}^n \to \mathbb{C}; ((x_i),(y_i)) \mapsto \Sigma(x_i y_i)$ |
| Kronecker product | $\mathrm{K}_{m \times n} : \mathbb{C}^m \times \mathbb{C}^n \to \mathbb{C}^{mn}; ((x_i),\mathbf{y})) \mapsto (x_i \mathbf{y})$ |
| *Others* | |
| Fork | $\mathrm{Fork}_n : \mathbb{C}^n \to \mathbb{C}^n \times \mathbb{C}^n; \mathbf{x} \mapsto (\mathbf{x},\mathbf{x})$ |
| Split | $\mathrm{Split}_n : \mathbb{C}^n \to \mathbb{C}^{n/2} \times \mathbb{C}^{n/2}; \mathbf{x} \mapsto (\mathbf{x}^U, \mathbf{x}^L)$ |
| Concatenate | $\oplus_n : \mathbb{C}^n \times \mathbb{C}^m \to \mathbb{C}^{n+m}; (\mathbf{x}, \mathbf{y}) \mapsto \mathbf{x} \oplus \mathbf{y}$ |
| Duplication | $\mathrm{dup}_n^m : \mathbb{C}^n \to \mathbb{C}^{nm}; (\mathbf{x} \mapsto \mathbf{x} \otimes I_m$ |
| Min | $\min_n : \mathbb{C}^n \times \mathbb{C}^n \to \mathbb{C}^n; (\mathbf{x}, \mathbf{y}) \mapsto (\min(x_i, y_i))$ |
| Max | $\max_n : \mathbb{C}^n \times \mathbb{C}^n \to \mathbb{C}^n; (\mathbf{x}, \mathbf{y}) \mapsto (\max(x_i, y_i))$ |

# OL Tensor Product: Repetitive Structure

**Kronecker product**

*(structured matrices)*

$$\mathrm{I}_m \otimes A_n = \begin{bmatrix} A_n & & \\ & \ddots & \\ & & A_n \end{bmatrix}$$

**OL Tensor product**

*(structured operators)*

$$\mathrm{I}_{M/M_B \times 1 \to M/M_B} \otimes \mathsf{MMM}_{M_B,N,K} = \left( (A, B) \mapsto \begin{bmatrix} A_0 \\ \vdots \\ A_{M/M_B - 1} \end{bmatrix} B \right)$$

**Definition**

*(extension to non-linear)*

$$(\mathrm{I}_{m \times n \to mn} \otimes \mathsf{A}) \left( \sum_{i=0}^{m-1} \mathbf{e}_i^m \otimes \mathbf{x}, \sum_{i=0}^{n-1} \mathbf{e}_i^n \otimes \mathbf{y} \right) = \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} \mathbf{e}_i^m \otimes \mathbf{e}_j^n \otimes \mathsf{A}(\mathbf{x}, \mathbf{y})$$

$$(\mathsf{A} \otimes \mathrm{I}_{m \times n \to mn}) \left( \sum_{i=0}^{m-1} \mathbf{x} \otimes \mathbf{e}_i^m, \sum_{i=0}^{n-1} \mathbf{y} \otimes \mathbf{e}_i^n \right) = \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} \mathsf{A}(\mathbf{x}, \mathbf{y}) \otimes \mathbf{e}_i^m \otimes \mathbf{e}_j^n$$

**SPIRAL**
www.spiral.net

# Translating OL Formulas Into Programs

| operator formula | code |
|---|---|
| $\mathbf{r} = (A_{K \times M \to N} \circ B_{k \times m \to K \times M})(\mathbf{x}, \mathbf{y})$ | ```(t, u) = B(x, y);```<br>```r = A(t, u);``` |
| $(\mathbf{r}, \mathbf{s}) = (A_{m \to M} \times B_{n \to N})(\mathbf{x}, \mathbf{y})$ | ```r = A(x);```<br>```s = B(y);``` |
| $\mathbf{r} = (I_{m \times n \to mn} \otimes A_{M \times N \to K})(\mathbf{x}, \mathbf{y})$ | ```for (i=0;i<m;i++)```<br>```    for (j=0;j<m;j++)```<br>```        r[(i*m+j)*K:1:(i*m+j+1)*K-1] =```<br>```            A(x[i*M:1:i*M+M-1], y[j*N:1:j*N+N-1]);``` |
| $\mathbf{r} = (A_{M \times N \to K} \otimes I_{m \times n \to mn})(\mathbf{x}, \mathbf{y})$ | ```for (i=0;i<m;i++)```<br>```    for (j=0;j<n;j++)```<br>```        r[i*m+j:m*n:i*m+j+m*n*(K-1))] =```<br>```            A(x[i:m:i+m*(M-1))], y[j:n:j+n*(N-1)]);``` |

$$I_{M/M_B \times 1 \to M/M_B} \otimes \mathsf{MMM}_{M_B, N, K} = \left( (A, B) \mapsto \left[ \frac{A_0}{\frac{\vdots}{A_{M/M_B - 1}}} \right] B \right)$$

# Example: Matrix Multiplication (MMM)

## Breakdown rules:

*capture various forms of blocking*

| breakdown rule | description |
| --- | --- |
| $\text{MMM}_{1,1,1} \rightarrow (\cdot)_1$ | base case |
| $\text{MMM}_{m,n,k} \rightarrow (\otimes)_{m/m_b \times 1} \otimes \text{MMM}_{m_b,n,k}$ | horizontal blocking  |
| $\text{MMM}_{m,n,k} \rightarrow \text{MMM}_{m,nb,k} \otimes (\otimes)_{1 \times n/nb}$ | interleaved blocking  |
| $\text{MMM}_{m,n,k} \rightarrow ((\Sigma_{k/k_b} \circ (\cdot)_{k/k_b}) \otimes \text{MMM}_{m,n,k_b}) \circ \\ ((L_{k/k_b}^{mk/k_b} \otimes I_{k_b}) \times I_{kn})$ | accumulative blocking  |
| $\text{MMM}_{m,n,k} \rightarrow (L_m^{mn/n_b} \otimes I_{n_b}) \circ \\ ((\otimes)_{1 \times n/n_b} \otimes \text{MMM}_{m,n_b,k}) \circ \\ (I_{km} \times (L_{n/n_b}^{kn/n_b} \otimes I_{n_b}))$ | vertical blocking  |

# Example: SAR Computation as OL Rules



$$\mathsf{SAR}_{k\times m\to n\times n} \;\to\; \mathsf{DFT}_{n\times n} \circ \mathsf{Interp}_{k\times m\to n\times n}$$

$$\mathsf{DFT}_{n\times n} \;\to\; (\mathsf{DFT}_n \otimes \mathrm{I}_n) \circ (\mathrm{I}_n \otimes \mathsf{DFT}_n)$$

$$\mathsf{Interp}_{k\times m\to n\times n} \;\to\; (\mathsf{Interp}_{k\to n} \otimes_i \mathrm{I}_n) \circ (\mathrm{I}_k \otimes_i \mathsf{Interp}_{m\to n})$$

$$\mathsf{Interp}_{r\to s} \;\to\; \left(\bigoplus_{i=0}^{n-2} \mathsf{InterpSeg}_k\right) \oplus \mathsf{InterpSegPruned}_{k,\ell}$$

$$\mathsf{InterpSeg}_k \;\to\; \mathrm{G}_f^{u\cdot n\to k} \circ \mathsf{iPrunedDFT}_{n\to u\cdot n} \circ \left(\frac{1}{n}\right) \circ \mathsf{DFT}_n$$

# Organization

- Operator language and algorithms

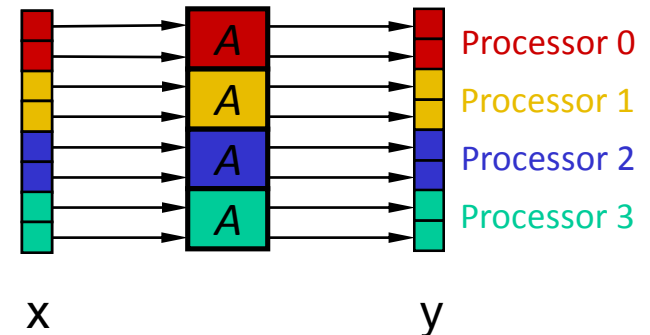- **Optimizing algorithms for platforms**

- Performance results

- Summary

# Modeling Multicore: Base Cases

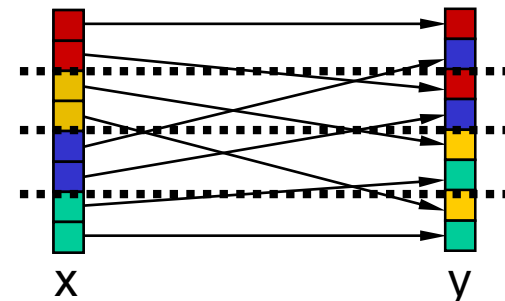- **Hardware abstraction: shared cache with cache lines**



$$\underbrace{A_{k \times m \to n}}_{\mathrm{smp}(p,\mu)}$$

- **Tensor product: embarrassingly parallel operator**

$$y = \Big( \mathbf{I}_p \otimes A \Big)(x)$$



- **Permutation: problematic; may produce false sharing**

$$y = \mathbf{L}_4^8(x)$$

# Parallelization: OL Rewriting Rules

- **Tags encode hardware constraints**
- **Rules are algorithm-independent**
- **Rules encode program transformations**

$$\underbrace{\left(\mathrm{I}_k \otimes \mathsf{L}_n^{mn}\right)}_{\mathsf{smp}(p,\mu)} \circ \underbrace{\mathsf{L}_{km}^{kmn}}_{\mathsf{smp}(p,\mu)} \to \left(\mathsf{L}_k^{kn} \otimes \mathrm{I}_{m/\mu}\right) \bar{\otimes}\, \mathrm{I}_\mu$$

$$\underbrace{\mathsf{L}_n^{kmn}}_{\mathsf{smp}(p,\mu)} \circ \underbrace{\left(\mathrm{I}_k \otimes \mathsf{L}_m^{mn}\right)}_{\mathsf{smp}(p,\mu)} \to \left(\mathsf{L}_n^{kn} \otimes \mathrm{I}_{m/\mu}\right) \bar{\otimes}\, \mathrm{I}_\mu$$

$$\underbrace{\mathsf{A} \circ \mathsf{B}}_{\mathsf{smp}(p,\mu)} \to \underbrace{\mathsf{A}}_{\mathsf{smp}(p,\mu)} \circ \underbrace{\mathsf{B}}_{\mathsf{smp}(p,\mu)}$$

$$\underbrace{\mathsf{A}^{k \times m \to n} \otimes \mathrm{I}^{1 \times p \to p}}_{\mathsf{smp}(p,\mu)} \to \underbrace{\mathsf{L}_n^{pn}}_{\mathsf{smp}(p,\mu)} \circ \left(\mathrm{I}_{1 \times p \to p} \otimes_{\parallel} \mathsf{A}^{k \times m \to n}\right) \circ \underbrace{\left(\mathrm{I}_k \times \mathsf{L}_p^{pm}\right)}_{\mathsf{smp}(p,\mu)}$$

$$\underbrace{\left(\mathsf{A} \times \mathsf{B}\right)}_{\mathsf{smp}(p,\mu)} \circ \underbrace{\left(\mathsf{C} \times \mathsf{D}\right)}_{\mathsf{smp}(p,\mu)} \to \underbrace{\left(\mathsf{A} \circ \mathsf{C}\right)}_{\mathsf{smp}(p,\mu)} \times \underbrace{\left(\mathsf{B} \circ \mathsf{D}\right)}_{\mathsf{smp}(p,\mu)}$$

# The Joint Rule Set: MMM

- ## Algorithm rules: breakdown rules

$$\text{MMM}_{m,n,k} \rightarrow (\otimes)_{m/m_b \times 1} \otimes \text{MMM}_{m_b,n,k}$$

$$\text{MMM}_{m,n,k} \rightarrow \text{MMM}_{m,nb,k} \otimes (\otimes)_{1 \times n/nb}$$

$$\text{MMM}_{m,n,k} \rightarrow ((\Sigma_{k/k_b} \circ (\cdot)_{k/k_b}) \otimes \text{MMM}_{m,n,k_b}) \circ$$
$$((L_{k/k_b}^{mk/k_b} \otimes I_{k_b}) \times I_{kn})$$

- ## Hardware constraints: base cases

$$\mathsf{I}_p \otimes_\parallel A_{m\mu \rightarrow n\mu}, \qquad \mathsf{P}_p \otimes_\parallel A_{k\mu \times m\mu \rightarrow n\mu}, \qquad \mathsf{K}_{q \times r} \otimes_\parallel A_{k\mu \times m\mu \rightarrow n\mu},$$
$$M \bar{\otimes} \mu, \qquad (M \bar{\otimes} \mathsf{I}_\mu) \otimes (N \bar{\otimes} \mathsf{I}_\mu)$$

- ## Program transformations: manipulation rules

$$\underbrace{\left( \mathsf{I}_k \otimes \mathsf{L}_n^{mn} \right)}_{\text{smp}(p,\mu)} \circ \underbrace{\mathsf{L}_{km}^{kmn}}_{\text{smp}(p,\mu)} \rightarrow \left( \mathsf{L}_k^{kn} \otimes \mathsf{I}_{m/\mu} \right) \bar{\otimes} \mathsf{I}_\mu$$

$$\underbrace{\mathsf{L}_n^{kmn}}_{\text{smp}(p,\mu)} \circ \underbrace{\left( \mathsf{I}_k \otimes \mathsf{L}_m^{mn} \right)}_{\text{smp}(p,\mu)} \rightarrow \left( \mathsf{L}_n^{kn} \otimes \mathsf{I}_{m/\mu} \right) \bar{\otimes} \mathsf{I}_\mu$$

*Combined rule set spans search space for empirical optimization*

# Parallelization Through Rewriting: MMM

$$\underbrace{\mathrm{MMM}_{m,n,k}}_{\mathrm{smp}(p,\mu)}$$

$$\rightarrow \underbrace{\left(\mathrm{I}_m \otimes \mathsf{L}_p^n\right) \circ \left(\mathrm{MMM}_{m,n/p,k} \otimes (\otimes)_{1\times p\rightarrow p}\right) \circ \left(\mathrm{I}_{km} \times (\mathrm{I}_k \otimes \mathsf{L}_{n/p}^n)\right)}_{\mathrm{smp}(p,\mu)}$$

$$\rightarrow \underbrace{\left(\mathrm{I}_m \otimes \mathsf{L}_p^n\right)}_{\mathrm{smp}(p,\mu)} \circ \underbrace{\left(\mathrm{MMM}_{m,n/p,k} \otimes (\otimes)_{1\times p\rightarrow p}\right)}_{\mathrm{smp}(p,\mu)} \circ \underbrace{\left(\mathrm{I}_{km} \times (\mathrm{I}_k \otimes \mathsf{L}_{n/p}^n)\right)}_{\mathrm{smp}(p,\mu)}$$

$$\rightarrow \underbrace{\left(\mathrm{I}_m \otimes \mathsf{L}_p^n\right)}_{\mathrm{smp}(p,\mu)} \circ \underbrace{\mathsf{L}_{m/pn}^{mn}}_{\mathrm{smp}(p,\mu)} \circ \left((\otimes)_{1\times p\rightarrow p} \otimes_{\parallel} \mathrm{MMM}_{m/p,n,k}\right) \circ \underbrace{\left(\mathrm{I}_{km} \times \mathsf{L}_p^{kn}\right)}_{\mathrm{smp}(p,\mu)} \circ \underbrace{\left(\mathrm{I}_{km} \times (\mathrm{I}_k \otimes \mathsf{L}_{n/p}^n)\right)}_{\mathrm{smp}(p,\mu)}$$

$$\rightarrow \left((\mathsf{L}_m^{mp} \otimes \mathrm{I}_{n/(p\mu)}) \bar{\otimes} \mathrm{I}_\mu\right) \circ \left((\otimes)_{1\times p\rightarrow p} \otimes_{\parallel} \mathrm{MMM}_{m,n/p,k}\right) \circ \left((\mathrm{I}_{km/\mu} \bar{\otimes} \mathrm{I}_\mu) \times ((\mathsf{L}_p^{kp} \otimes \mathrm{I}_{n/(p\mu)}) \bar{\otimes} \mathrm{I}_\mu)\right)$$

**Load-balanced**
**No false sharing**

# Same Approach for Different Paradigms

### Threading:

$$\underset{\mathsf{smp}(p,\mu)}{\underline{\mathbf{DFT}_{mn}}} \;\rightarrow\; \underbrace{\left((\mathbf{DFT}_m \otimes \mathbf{I}_n)\mathsf{T}_n^{mn}(\mathbf{I}_m \otimes \mathbf{DFT}_n)\mathsf{L}_m^{mn}\right)}_{\mathsf{smp}(p,\mu)}$$

$$\cdots$$

$$\rightarrow\; \underset{\mathsf{smp}(p,\mu)}{\underline{\left(\mathbf{DFT}_m \otimes \mathbf{I}_n\right)}} \; \underset{\mathsf{smp}(p,\mu)}{\underline{\mathsf{T}_n^{mn}}} \; \underset{\mathsf{smp}(p,\mu)}{\underline{\left(\mathbf{I}_m \otimes \mathbf{DFT}_n\right)}} \; \underset{\mathsf{smp}(p,\mu)}{\underline{\mathsf{L}_m^{nm}}}$$

$$\cdots$$

$$\rightarrow\; \left((\mathsf{L}_m^{mp} \otimes \mathbf{I}_{n/p\mu}) \otimes_\mu \mathbf{I}_\mu\right)\left(\mathbf{I}_p \otimes_\| (\mathbf{DFT}_m \otimes \mathbf{I}_{n/p})\right)\left((\mathsf{L}_p^{mp} \otimes \mathbf{I}_{n/p\mu}) \otimes_\mu \mathbf{I}_\mu\right)$$

$$\left(\bigoplus_{i=0}^{p-1}{}_\| \mathsf{T}_n^{mn,i}\right)\left(\mathbf{I}_p \otimes_\|(\mathbf{I}_{m/p} \otimes \mathbf{DFT}_n)\right)\left(\mathbf{I}_p \otimes_\| \mathsf{L}_{m/p}^{mn/p}\right)\left((\mathsf{L}_p^{pn} \otimes \mathbf{I}_{m/p\mu}) \otimes_\mu \mathbf{I}_\mu\right)$$

### Vectorization:

$$\underset{\mathsf{vec}(\nu)}{\underline{(\overline{\mathbf{DFT}_{mn}})}} \;\rightarrow\; \underbrace{\left((\mathbf{DFT}_m \otimes \mathbf{I}_n)\mathsf{T}_n^{mn}(\mathbf{I}_m \otimes \mathbf{DFT}_n)\mathsf{L}_m^{mn}\right)}_{\mathsf{vec}(\nu)}$$

$$\cdots$$

$$\rightarrow\; \underset{\mathsf{vec}(\nu)}{\underline{(\overline{\mathbf{DFT}_m \otimes \mathbf{I}_n})^\nu}} \underset{\mathsf{vec}(\nu)}{\underline{(\overline{\mathsf{T}_n^{mn}})^\nu}} \underset{\mathsf{vec}(\nu)}{\underline{(\overline{\mathbf{I}_m \otimes \mathbf{DFT}_n})\mathsf{L}_m^{mn}}^\nu}$$

$$\cdots$$

$$\rightarrow\; (\mathbf{I}_{mn/\nu} \otimes \underset{\mathsf{sse}}{\underline{\mathsf{L}_\nu^{2\nu}}})(\overline{\mathbf{DFT}_m \otimes \mathbf{I}_{n/\nu}} \vec{\otimes} \mathbf{I}_\nu)(\overline{\mathsf{T}_n^{mn}})^\nu$$

$$\left(\mathbf{I}_{m/\nu} \otimes (\overline{\mathsf{L}_\nu^n} \vec{\otimes} \mathbf{I}_\nu)(\mathbf{I}_{n/\nu} \otimes (\mathsf{L}_\nu^{2\nu} \vec{\otimes} \mathbf{I}_\nu)(\mathbf{I}_2 \otimes \underset{\mathsf{sse}}{\underline{\mathsf{L}_\nu^{\nu^2}}})(\mathsf{L}_2^{2\nu} \vec{\otimes} \mathbf{I}_\nu))(\overline{\mathbf{DFT}_n} \vec{\otimes} \mathbf{I}_\nu)\right)$$

$$\left((\mathsf{L}_m^{mn} \otimes \mathbf{I}_2) \vec{\otimes} \mathbf{I}_\nu\right)(\mathbf{I}_{mn/\nu} \otimes \underset{\mathsf{sse}}{\underline{\mathsf{L}_2^{2\nu}}})$$

### GPUs:

$$\underset{\mathsf{gpu}(t,c)}{\underline{\left(\mathbf{DFT}_{r^k}\right)}} \;\rightarrow\; \underbrace{\left(\prod_{i=0}^{k-1}\mathsf{L}_r^{r^k}\left(\mathbf{I}_{r^{k-1}} \otimes \mathbf{DFT}_r\right)\left(\mathsf{L}_{r^{k-i-1}}^{r^k}(\mathbf{I}_{r^i} \otimes \mathsf{T}_{r^{k-i-1}}^{r^{k-i}})\underset{\mathsf{vec}(c)}{\underline{\mathsf{L}_{r^{i+1}}^{r^k}}}\right)\right)\mathsf{R}_r^{r^k}}_{\mathsf{gpu}(t,c)}$$

$$\cdots$$

$$\rightarrow\; \left(\prod_{i=0}^{k-1}(\mathsf{L}_r^{r^n/2} \vec{\otimes} \mathbf{I}_2)\left(\mathbf{I}_{r^{n-1}/2} \otimes_\times \underset{\mathsf{shd}(t,c)}{\underline{(\mathbf{DFT}_r \vec{\otimes} \mathbf{I}_2)\mathsf{L}_r^{2r}}}\right)\mathsf{T}_i\right)$$

$$(\mathsf{L}_r^{r^n/2} \vec{\otimes} \mathbf{I}_2)(\mathbf{I}_{r^{n-1}/2} \otimes_\times \underset{\mathsf{shd}(t,c)}{\underline{\mathsf{L}_r^{2r}}})(\mathsf{R}_r^{r^{n-1}} \vec{\otimes} \mathbf{I}_r)$$

### Verilog for FPGAs:

$$\underset{\mathsf{stream}(r^s)}{\underline{\left(\mathbf{DFT}_{r^k}\right)}} \;\rightarrow\; \underbrace{\left[\prod_{i=0}^{k-1}\mathsf{L}_r^{r^k}\left(\mathbf{I}_{r^{k-1}} \otimes \mathbf{DFT}_r\right)\left(\mathsf{L}_{r^{k-i-1}}^{r^k}(\mathbf{I}_{r^i} \otimes \mathsf{T}_{r^{k-i-1}}^{r^{k-i}})\mathsf{L}_{r^{i+1}}^{r^k}\right)\right]\mathsf{R}_r^{r^k}}_{\mathsf{stream}(r^s)}$$

$$\cdots$$

$$\rightarrow\; \left[\prod_{i=0}^{k-1}\underset{\mathsf{stream}(r^s)}{\underline{\mathsf{L}_r^{r^k}}}\underset{\mathsf{stream}(r^s)}{\underline{\left(\mathbf{I}_{r^{k-1}} \otimes \mathbf{DFT}_r\right)}}\underset{\mathsf{stream}(r^s)}{\underline{\left(\mathsf{L}_{r^{k-i-1}}^{r^k}(\mathbf{I}_{r^i} \otimes \mathsf{T}_{r^{k-i-1}}^{r^{k-i}})\mathsf{L}_{r^{i+1}}^{r^k}\right)}}\right]\underset{\mathsf{stream}(r^s)}{\underline{\mathsf{R}_r^{r^k}}}$$

$$\cdots$$

$$\rightarrow\; \left[\prod_{i=0}^{k-1}\underset{\mathsf{stream}(r^s)}{\underline{\mathsf{L}_r^{r^k}}}\underset{\mathsf{stream}(r^s)}{\underline{\left(\mathbf{I}_{r^{k-s-1}} \otimes_s (\mathbf{I}_{r^{s-1}} \otimes \mathbf{DFT}_r)\right)}}\underset{\mathsf{stream}(r^s)}{\underline{\mathsf{T}_i'}}\right]\underset{\mathsf{stream}(r^s)}{\underline{\mathsf{R}_r^{r^k}}}$$
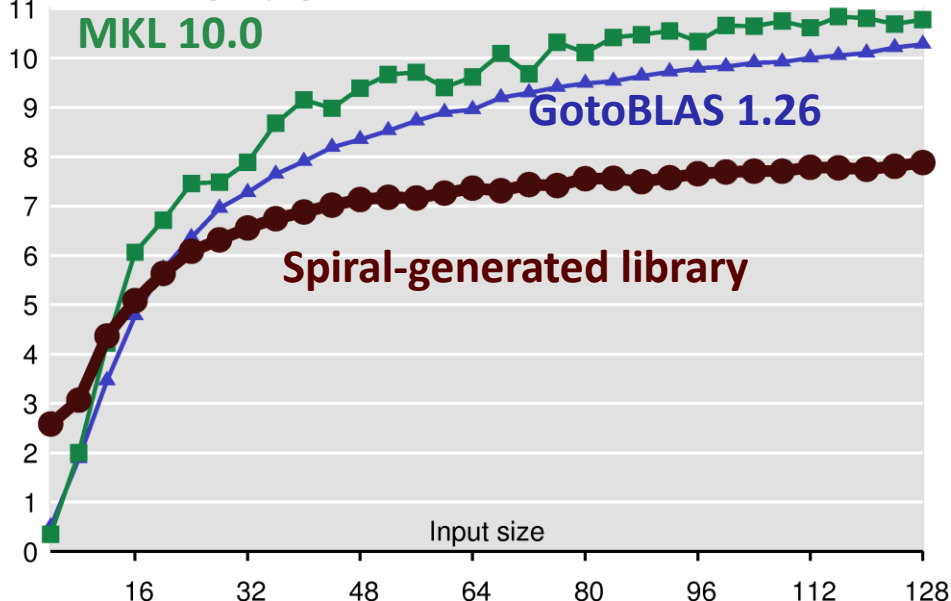
# Organization

- **Operator language and algorithms**

- **Optimizing algorithms for platforms**

- **Performance results**
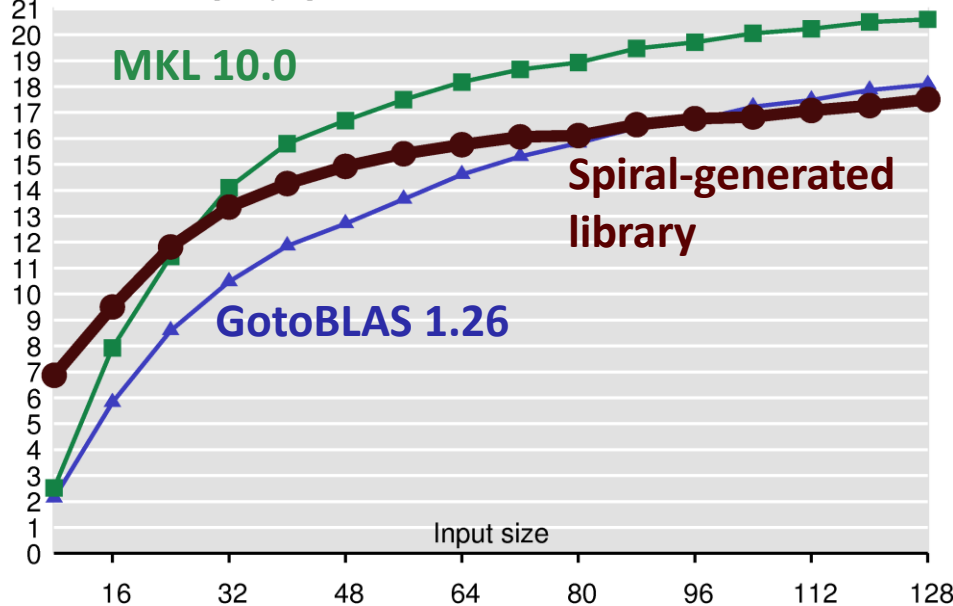
- **Summary**

# Matrix Multiplication Library

**GEMM (acc. packed sq. NN), single-threaded, double precision**

Performance [Gflop/s]                    Dual Intel Xeon 5160, 3000 MHz, icc 10.1
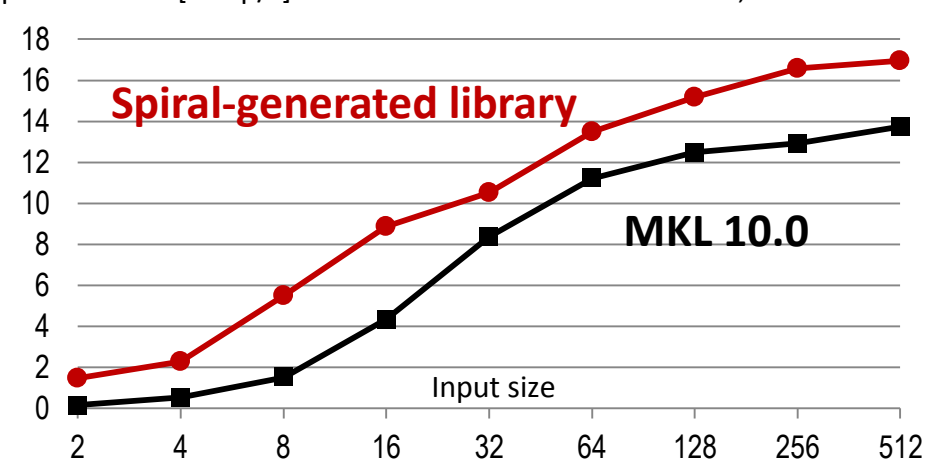


MKL 10.0

GotoBLAS 1.26

Spiral-generated library

Input size

**GEMM (acc. packed sq. NN), single-threaded, single precision**

Performance [Gflop/s]                    Dual Intel Xeon 5160, 3000 MHz, icc 10.1



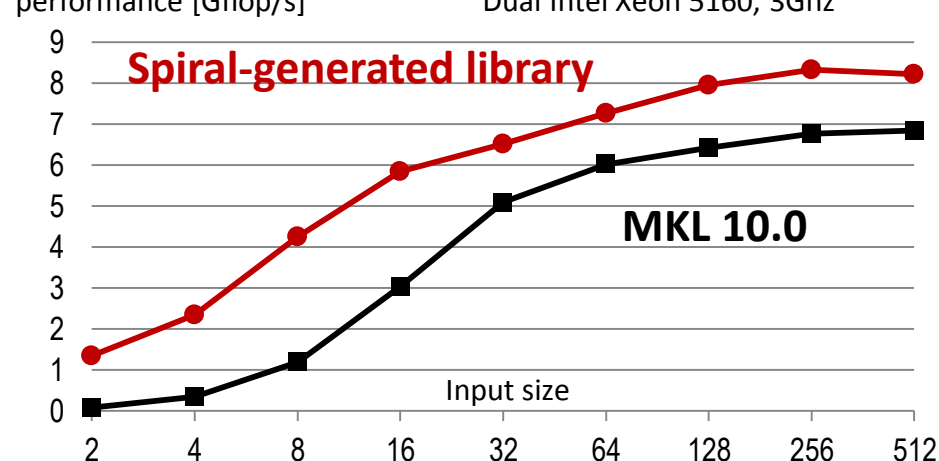MKL 10.0

Spiral-generated library

GotoBLAS 1.26

Input size

**Rank-k Update**, single precision, k=4

performance [Gflop/s]                    Dual Intel Xeon 5160, 3Ghz



Spiral-generated library

MKL 10.0

Input size

**Rank-k Update**, double precision, k=4

performance [Gflop/s]                    Dual Intel Xeon 5160, 3Ghz



Spiral-generated library

MKL 10.0

Input size

**SPIRAL**
www.spiral.net

# Result: Spiral-Generated PFA SAR on Core2 Quad

**SAR Image Formation on Intel platforms**

performance [Gflop/s]



Legend:
- 3.0 GHz Core 2 (65nm)
- 3.0 GHz Core 2 (45nm)
- 2.66 GHz Core i7
- 3.0 GHz Core i7 (Virtual)

*newer platforms*

- **Algorithm by J. Rudin (best paper award, HPEC 2007): 30 Gflop/s on Cell**
- **Each implementation: vectorized, threaded, cache tuned, ~13 MB of code**

# Organization

- **Operator language and algorithms**

- **Optimizing algorithms for platforms**

- **Performance results**

- **Summary**

# Summary

- **Platforms are powerful yet complicated**
  optimization will stay a hard problem


Image: Intel

- **OL: unified mathematical framework**
  captures platforms and algorithms

$$\underbrace{\mathrm{I}_p \otimes A_n}_{\mathsf{smp}(p,\mu)}$$

- **Spiral: program generation and autotuning**
  can provide full automation



- **Performance of supported kernels**
  is competitive with expert tuning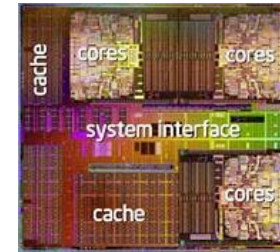